



Particular Object Retrieval With Integral Max-Pooling of CNN Activations

Giorgos Tolias, Ronan Sire, Hervé Jégou

► To cite this version:

Giorgos Tolias, Ronan Sire, Hervé Jégou. Particular Object Retrieval With Integral Max-Pooling of CNN Activations. ICLR 2016 - International Conference on Learning Representations, May 2016, San Juan, Puerto Rico. pp.1-12. hal-01842218

HAL Id: hal-01842218

<https://inria.hal.science/hal-01842218>

Submitted on 18 Jul 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

PARTICULAR OBJECT RETRIEVAL WITH INTEGRAL MAX-POOLING OF CNN ACTIVATIONS

Giorgos Tolias *
Center for Machine Perception
FEE CTU Prague

Ronan Sifre
Irisa Rennes

Hervé Jégou
Facebook AI Research

ABSTRACT

Recently, image representation built upon Convolutional Neural Network (CNN) has been shown to provide effective descriptors for image search, outperforming pre-CNN features as short-vector representations. Yet such models are not compatible with geometry-aware re-ranking methods and still outperformed, on some particular object retrieval benchmarks, by traditional image search systems relying on precise descriptor matching, geometric re-ranking, or query expansion. This work revisits both retrieval stages, namely initial search and re-ranking, by employing the same primitive information derived from the CNN. We build compact feature vectors that encode several image regions without the need to feed multiple inputs to the network. Furthermore, we extend integral images to handle max-pooling on convolutional layer activations, allowing us to efficiently localize matching objects. The resulting bounding box is finally used for image re-ranking. As a result, this paper significantly improves existing CNN-based recognition pipeline: We report for the first time results competing with traditional methods on the challenging Oxford5k and Paris6k datasets.

1 INTRODUCTION

CONTENT based image retrieval has received a sustained attention over the last decade, leading to mature systems for tasks like visual instance retrieval. Current state-of-the-art approaches are derived from the Bag-of-Words model of Sivic & Zisserman (2003) and mainly owe their success to locally invariant features (Lowe, 2004) and large visual codebooks (Philbin et al., 2007). These methods are typically composed of an initial *filtering* stage where all database images are ranked in terms of similarity to a query image and a second *re-ranking* stage, which refines the search results of the top-ranked elements. The filtering stage is improved in several ways, such as incorporating weak geometric information (Jégou et al., 2010), employing compact approximations of the local descriptors (Jégou et al., 2010), or learning smart codebooks (Mikulik et al., 2013; Avrithis & Kalantidis, 2012). In such cases, local descriptors are individually matched and selective matching functions (Tolias et al., 2015; Tao et al., 2014) improve the search quality. Geometric matching models (Philbin et al., 2007; Avrithis & Tolias, 2014) are typically applied in a pairwise manner during the re-ranking stage of a short-list of images. Query expansion approaches significantly increase the performance (Chum et al., 2011), at the cost of larger query times.

The recent advances achieved by Convolutional Neural Networks (CNN) and the use of intermediate layer activations as feature vectors (Donahue et al., 2013) create opportunities for representations that are competitive for image or particular object retrieval, and not only classification tasks. Several works have already investigated this research direction, such as global or local representations based on either fully connected (Babenko et al., 2014; Gong et al., 2014) or convolutional layers (Razavian et al., 2014b; Azizpour et al., 2014; Babenko & Lempitsky, 2015). The performance of CNN-based features has rapidly improved to the point of competing and even outperforming pre-CNN works that aggregate local features (Jégou et al., 2012; Radenović et al., 2015). In particular, activations of convolutional layers followed by a global max-pooling operation (Azizpour et al., 2014) produce highly competitive compact image representations. One limitation is that such approaches are not compatible with the geometric-aware models involved in the final re-ranking stages.

*Research partially conducted while G. Tolias and H. Jégou were at Inria. We would like to thank Florent Perronnin for his valuable feedback. This work was partly supported by MSMT LL1303 ERC-CZ grant.

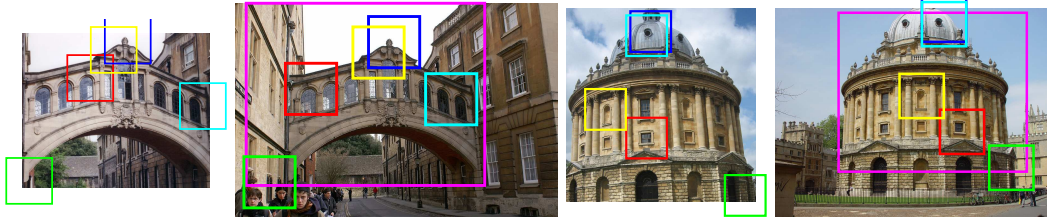


Figure 1: Query objects (left) and the corresponding localization in another image (right) are shown. We visualize the patches that contribute the highest to the image similarity score. Displayed patches correspond to the receptive field of CNN activations. Object localization is displayed in magenta, while different colors are used for patches in correspondence.

This work revisits both filtering and re-ranking stages with CNN-based features. We make the three following contributions.

- First, we propose a compact image representation derived from the convolutional layer activations that encodes multiple image regions without the need to re-feed multiple inputs to the network, in spirit of recent Fast-RCNN (Girshick, 2015) and Faster-RCNN (Ren et al., 2015) methods but here targeting particular object retrieval. The underlying primitive representation is used in all stages (initial retrieval and re-ranking).
- Second, we employ the generalized mean (Dollár et al., 2009) to enable the use of integral images along with max-pooling. This efficient method is exploited for particular object localization (see Figure 1) directly in the 2D maps of CNN activations.
- Third, our localization approach is used for image re-ranking and leads us to define a simple yet effective query expansion method.

These approaches are complementary and, when combined, produce for the first time a system which compete on the Oxford and Paris building benchmarks with state-of-the-art re-ranking approaches based on local features. Our approach outperforms by a large margin previous methods based on CNN, while being more efficient in practice.

2 RELATED WORK

CNN based representation. A typical CNN consists of several convolutional layers, followed by fully connected layers and ends with a softmax layer producing a distribution over the training classes. Instead of using this inherent classifier, one can consider the activations of the intermediate layers to train a classifier. In particular, the activations of the fully connected layers have been shown to be very effective and capable of adaptation to various domains (Oquab et al., 2014), such as scene recognition (Donahue et al., 2013; Sicre & Jurie, 2015), object detection (Iandola et al., 2014), and semantic segmentation (Girshick et al., 2014). In the case of image retrieval, fully connected layers are used as global descriptors followed by dimensionality reduction (Babenko et al., 2014). They are also employed as region descriptors to be compared to database descriptors (Razavian et al., 2014a) or aggregated in a VLAD manner (Gong et al., 2014).

Recent works derive visual representations from the activations of the convolutional layers. This is achieved either by stacking activations (Girshick et al., 2014) or by performing spatial max-pooling (Azizpour et al., 2014) or sum-pooling (Babenko & Lempitsky, 2015) for each feature channel. According to Azizpour et al. (2014) such representation offers better generalization properties for test data that are far from the source (training) data. Noticeably, higher performance in particular object or scene retrieval is obtained by using convolutional layers rather than fully connected ones. The very recent work of Babenko & Lempitsky (2015) shows that sum-pooling performs better than max-pooling when the image representation is whitened. In addition to be a costly choice, we will show that this is not optimal in our context of object localization (see Section 8). Finally, Kalantidis et al. (2015) propose spatial and feature channel weighting that significantly improves performance. Their approach is complementary to what we propose for the filtering and the re-ranking stage.

Recent examples utilize information from fully connected layers to perform generic object detection (Iandola et al., 2014; Papandreou et al., 2014). Such approaches are prohibitive for the re-ranking purposes of large scale image retrieval. They have high computational cost and the inherent features are not optimal for particular object matching.

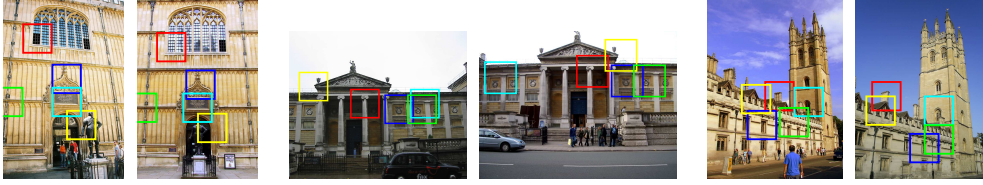


Figure 2: We visualize the receptive fields related to the 5 MAC components that contribute the most to the image similarity. Each displayed receptive field corresponds to the maximum response of a feature channel. A different color is used for each feature channel, while different feature channels are shown for each image pair.

Localization. In the recent years, the sliding window principle has been quite successful for many object localization methods. Due to the large number of possible windows, exhaustive search is extremely costly. However, integral images (Viola & Jones, 2001) offer a constant cost solution to the evaluation of a single region. This attractive alternative is applicable for feature vectors constructed via a sum-pooling operation.

A globally optimal solution is given by Efficient Subwindow Search (ESS) of Lampert et al. (2009), who use branch-and-bound search to avoid exhaustive search. Their work employs integral images, which are also used in later improvements of ESS (An et al., 2009). An et al. (2009) formalize localization as a maximum sub-array problem and similarly to Chen et al. (2013) they employ Bentley’s algorithm (Bentley, 1999). Integral images facilitate the evaluation of many region candidates (Uijlings et al., 2013) based on VLAD or Fisher vectors (Van de Sande et al., 2014). All aforementioned approaches take advantage of integral images due to the inherent sum-pooling operation in the given representation. In this paper, we extend integral images to perform max-pooling over CNN activation maps, which is shown to be a better choice for describing regions (as opposed to the entire image).

Several object localization techniques have been proposed in the context of image retrieval as well. Lampert (2009) propose a two layer branch-and-bound method that alternates between regions and images. Integral images offer a significant speed-up in the work of Lin & Brandt (2010) to perform localization through Bag-of-Words. The overall idea bears similarities with our work. However, we differentiate by employing CNN-based representation with max-pooling. Some approaches (Tao et al., 2014; Shen et al., 2014) individually index local features for localization. In our case, the localization method is built on top of a compact representation, initially used for the filtering stage. Finally, Arandjelovic & Zisserman (2013) propose a localization strategy based on VLAD, where similarity is computed for multiple image regions, giving a more precise localization via regression.

3 BACKGROUND

We consider a pre-trained CNN and discard all the fully connected layers. Given an input image I of size $W_I \times H_I$, the activations (responses) of a convolutional layer form a 3D tensor of $W \times H \times K$ dimensions, where K is the number of output feature channels, i.e. multi-dimensional filters. The spatial resolution $W \times H$ depends on the network architecture, the layer examined, and the input image resolution. We assume that Rectified Linear Units (ReLU) are applied as a last step, guaranteeing that all elements are non-negative.

We represent this 3D tensor of responses as a set of 2D feature channel responses $\mathcal{X} = \{\mathcal{X}_i\}$, $i = 1 \dots K$, where \mathcal{X}_i is the 2D tensor representing the responses of the i^{th} feature channel over the set Ω of valid spatial locations, and $\mathcal{X}_i(p)$ is the response at a particular position p . Therefore, the feature vector constructed by a spatial max-pooling over all locations (Azizpour et al., 2014) is given by

$$\mathbf{f}_\Omega = [f_{\Omega,1} \dots f_{\Omega,i} \dots f_{\Omega,K}]^\top, \text{ with } f_{\Omega,i} = \max_{p \in \Omega} \mathcal{X}_i(p). \quad (1)$$

Maximum activations of convolutions (MAC). Two images are compared with the cosine similarity of the K -dimensional vectors produced as described above. This representation, referred to as MAC, does not encode the location of the activations (unlike activations of fully connected layers),

due to the max-pooling operated over a single region of size $W \times H$. It encodes the maximum “local” response of each of the convolutional filters and is therefore translation invariant. In all the following, we consider the last convolutional layer of the examined networks.

Figure 2 visualizes the patches that contribute the most to the image similarity. They correspond either to the same object part or similar parts due to repeated structures. We extract MAC from input images of any resolution or aspect ratio by simply subtracting the mean pixel value (Iandola et al., 2014) from the input images. No crop or change of aspect ratio is required (Azizpour et al., 2014).

The max pooling operation that is performed over a single cell offers translation invariance to the resulting representation. This is in contrast to representation derived from the fully connected layers that requires objects to be aligned. In our case, we assume that objects are up-right and we simply benefit from the rotation tolerance provided by the CNN due to the training data used. The same stands for the tolerance to scale changes.

4 ENCODING REGIONS INTO SHORT VECTORS

This section describes how we exploit the activations of the CNN convolutional layers to derive representations for image regions. Region vectors are aggregated to produce a short signature used in the filtering stage of image retrieval.

Region feature vector. The feature vector \mathbf{f}_Ω described in Section 3 is a representation for the whole image I . Now, we consider a rectangular region $\mathcal{R} \subseteq \Omega = [1, W] \times [1, H]$, and define the regional feature vector

$$\mathbf{f}_\mathcal{R} = [\mathbf{f}_{\mathcal{R},1} \dots \mathbf{f}_{\mathcal{R},i} \dots \mathbf{f}_{\mathcal{R},K}]^\top \quad (2)$$

where $\mathbf{f}_{\mathcal{R},i} = \max_{p \in \mathcal{R}} \mathcal{X}_i(p)$ is the maximum activation of the i^{th} channel on the considered region.

The regions \mathcal{R} are defined on the space Ω of all valid positions for the considered feature map (and not on the input image plane). A region of size 1 corresponds to a feature vector consisting of a single activation at a particular location. We are now able to construct a representation for multiple regions without re-feeding additional input to the CNN, similarly to recent RNN variants (Ren et al., 2015; Girshick, 2015), which drastically reduces the processing cost.

Now assume a linear mapping of a given region \mathcal{R} back to the original image. The proposed region vector captures a larger image region than the back-projected one, due to the large receptive field. A similar effect occurs in the context of object detection (Iandola et al., 2014), where fully connected layers are applied in a sliding window fashion.

R-MAC: regional maximum activation of convolutions. We now consider a set of R regions of different sizes. The structure of the regions is similar to the one proposed by Razavian et al. (2014b), but we define them on the CNN response maps and not on the original image. We sample square regions at L different scales. At the largest scale ($l = 1$), the region size is determined to be as large as possible, i.e., its height and width are both equal to $\min(W, H)$. The regions are sampled uniformly such that the overlap between consecutive regions is as close as possible to 40%. Remark that the aspect ratio of the original image has an influence on the number m of regions that we extract (1 region only if the input image is square). At every other scale l we uniformly sample $l \times (l + m - 1)$ regions of width $2 \min(W, H)/(l + 1)$, as illustrated in Figure 3 (left).

Then we calculate the feature vector associated with each region, and post-process it with ℓ_2 -normalization, PCA-whitening (Jégou & Chum, 2012) and ℓ_2 -normalization. We combine the collection of regional feature vectors into a single image vector by summing them and ℓ_2 -normalizing in the end. This choice keeps the dimensionality low which is equal to the number of feature channels. However, we show in our experiments that the resulting representation, referred to as R-MAC, offers a significant better performance than the corresponding MAC with same dimensionality. Note, the aggregation of the region vectors can be seen as a simple kernel that cross matches all possible regions, including across different scale.

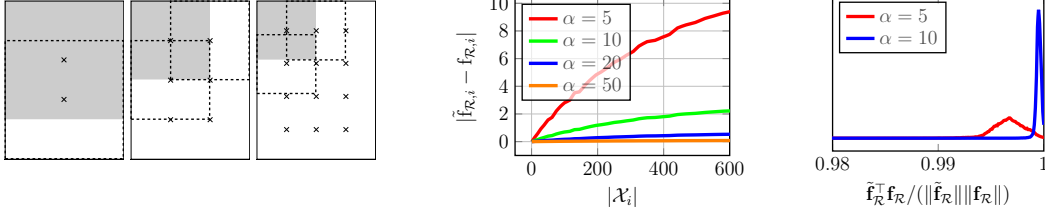


Figure 3: Left: Sample regions extracted at 3 different scales ($l = 1 \dots 3$). We show the top-left region of each scale (gray colored region) and its neighboring regions towards each direction (dashed borders). We depict the centers of all regions with a cross. Middle: Approximation error of the maximum value versus the size of the response set for different values of exponent α . Measurements are performed on 10 randomly selected images by evaluating all possible regions. The responses for this set of images take values in $[0, 151]$. Right: Empirical distribution of the cosine similarity value between the exact vector $\mathbf{f}_{\mathcal{R}}$ and its approximation $\tilde{\mathbf{f}}_{\mathcal{R}}$. Measurements are collected by constructing the exact and approximate vectors of all possible regions on 10 randomly sampled images.

5 OBJECT LOCALIZATION

In this section we propose an extension of integral images to perform approximate max-pooling over a set \mathcal{X} of 2D feature channel response maps, which provide a rough yet efficient localization to our CNN-based method.

Approximate integral max-pooling. Noticing that the responses \mathcal{X}_i are non-negative, we exploit the generalized mean (Dollár et al., 2009) to approximate each feature value $f_{\mathcal{R},i}$ associated with a given region \mathcal{R} by the estimate

$$\tilde{f}_{\mathcal{R},i} = \left(\sum_{p \in \mathcal{R}} \mathcal{X}_i(p)^\alpha \right)^{\frac{1}{\alpha}} \approx \max_{p \in \mathcal{R}} \mathcal{X}_i(p) = f_{\mathcal{R},i}, \quad (3)$$

where the parameter $\alpha > 1$ is such that $\tilde{f}_i \rightarrow f_i$ when $\alpha \rightarrow +\infty$.

Figure 3 (middle) shows the average approximation error $|\tilde{f}_{\mathcal{R},i} - f_{\mathcal{R},i}|$ estimated over several image regions. We report the approximation error as a function of the size of the corresponding response set on which the maximum value is computed. The various sizes of response sets are an outcome of using all possible regions. A high value of the exponent α leads to a better approximation, while applying on more elements makes the approximation less precise.

By approximating the maximum in this manner, we can now use integral images (Viola & Jones, 2001) to approximate the regional feature vector $\mathbf{f}_{\mathcal{R}}$ defined on any rectangular region \mathcal{R} . For each channel, we construct the integral image of the 2D tensor whose value at position p is equal to $\mathcal{X}_i(p)^\alpha$, $p \in \mathcal{R}$. Then, the sum of Equation (3) is simply given by the sum of 4 terms (Viola & Jones, 2001). This allow us to efficiently compute max-pooling for many regions and therefore to construct the corresponding feature vectors. This is in contrast to the explicit construction of many regions with representation derived from fully connected layers, which is prohibitive due to the need to resize/crop and re-feed each region to the network.

We evaluate the approximation quality by measuring the cosine similarity between the exact vector and its approximate counterpart. The distribution of this similarity is presented in Figure 3 (right) and is measured on all possible regions of 10 randomly selected images. The proposed approximation is very precise even for moderate values of α . We set α equal to 10 in all of our experiments.

Window detection. Let us now assume that there is another image Q depicting a single object, *i.e.* cropped via a bounding box defining the object of interest. We denote as \mathbf{q} the corresponding MAC feature vector. The 2D region, defined on the CNN activations \mathcal{X} of image I , that maximizes the similarity to \mathbf{q} is computed as

$$\hat{\mathcal{R}} = \arg \max_{\mathcal{R} \subseteq \Omega} \frac{\tilde{\mathbf{f}}_{\mathcal{R}}^T \mathbf{q}}{\|\tilde{\mathbf{f}}_{\mathcal{R}}\| \|\mathbf{q}\|}. \quad (4)$$

The region $\hat{\mathcal{R}}$ maximizing the similarity is mapped back to the original image I with a precision of $(\frac{W}{W_I}, \frac{H}{H_I})$ pixels, providing a rough localization of the object depicted in Q . The corresponding similarity does not take into account all the visual content of image I and is therefore free from the influence of background clutter. The brute-force detection of the optimal region by exhaustive search is expensive, as the number of possible regions is in $\mathcal{O}(W^2H^2)$. In preliminary tests, we have evaluated a globally optimal solution based on *branch and bound* search, as in ESS (Lampert et al., 2009). The necessary bounds are trivially derived for our representation. The search is not significantly sped up in our case: The maxima are not distinct enough and a large number of regions are considered, while the overhead of maintaining the priority queue is high.

AML: approximate max-pooling localization. Instead, we restrict the number of regions that we evaluate and locally refine the best ones with simple heuristics. Candidate regions are uniformly sampled with a *search step* equal to t . In addition, regions having an aspect ratio larger than s times that of the query region are discarded. The parameters of the best region are refined in a coordinate descent manner, while allowing a maximum change of 3 units. The refinement process is repeated up to 5 times. Experiments show that the overlap of the detected region to the optimal one is high.

6 RETRIEVAL, LOCALIZATION AND RE-RANKING

Initial retrieval. The MAC or R-MAC feature vector is computed for all databases images. Similarly, at query time we process the query image and extract the corresponding feature vector. During the filtering stage we directly evaluate cosine similarity between the query and all the database vectors. Therefore, we obtain the initial ranking based on the similarity of MAC or R-MAC vectors.

Re-ranking. We consider a second re-ranking stage, as typically performed in spatial verification (Philbin et al., 2007) with local features. A short-list of N top-ranked images is considered and AML, as described in Section 5, is applied on pairs of query and database images. Note that the query is now represented by the MAC vector, since this is used in AML, while the database image is represented by \mathcal{X} . For each re-ranked image we obtain a score given by the region that maximizes the similarity to the query. This similarity is used to re-rank the elements of the short-list. Furthermore, a rough localization of the query object is available.

Remarks: At the filtering stage, whitened MAC (whitening as described in Section 8) or R-MAC can be used, while the localization procedure employs similarity with respect to ℓ_2 -normalized MAC. However, once the query object is localized, then, similarity between the query and the detected region is computed via whitened MAC or R-MAC, depending on the chosen filtering method. This similarity score is used to perform re-ranking. The required representation is constructed on query time only for the detected region and is acquired efficiently with integral images.

Query expansion (QE). Re-ranking brings positive images at the very top ranked positions. Then, we collect the 5 top-ranked images, merge them with the query vector, and compute their mean. Finally, the similarity to this mean vector is adopted to re-rank once more the top N images.

7 IMPLEMENTATION DETAILS

We observe that thresholding the response values of \mathcal{X} which are larger than 128 (0.001% of all responses) and mapping each value to the closest smaller integer (floor operation) leads to insignificant losses. This allows the computation of α -th power with a lookup table and speeds-up the construction of integral images. Moreover, we approximate the α -th root of Equation (3) by performing binary search on the same lookup table of α -th power. This process allows the optimal window search to be more efficient.

The response maps represented by \mathcal{X} are sparse (Agrawal et al., 2014). In particular, using the network of Krizhevsky et al. (2012) on Oxford buildings dataset (Philbin et al., 2007) results in 81% of response values being zero, which is convenient for storage purpose. We further decrease the memory requirements by uniformly quantizing the responses into 8 values. This results in more elements mapped to the same value. Therefore, we store the positions of non-zeros values with delta coding and use only 1 byte per non-zero element. Note that an image of resolution equal to 1024×768 corresponds to feature channel response maps of size 30×22 using the same network. Finally, an image requires around 32 kB of memory. At re-ranking time we construct one integral image at a time and use double precision (8 bytes) for its elements.

Table 1: Left: Comparison between the exhaustive sliding window and our alternative of window sampling and refinement. We report the average IoU *w.r.t.* the globally optimal window and the average percentage of windows evaluated *w.r.t.* to the exhaustive search (noted by %W). Measurements are conducted on all pairs of Oxford5k query images and their corresponding positive images. Right: Performance (mAP) of MAC and R-MAC on Oxford5k. Resol. corresponds to the input image resolution (maximum dimension).

Search step t	Aspect ratio change threshold s					
	1.1		1.5		2.0	
	IoU	%W	IoU	%W	IoU	%W
1	81.8	8.9	88.7	27.5	93.7	46.3
2	79.9	0.5	83.8	2.0	86.6	3.6
3	78.7	0.2	81.2	0.5	83.6	0.8
4	77.0	0.1	79.5	0.2	81.5	0.3
5	75.8	0.1	79.0	0.1	80.9	0.2

Network	Resol.	MAC	R-MAC			
			$L=1$	$L=2$	$L=3$	$L=4$
AlexNet	1024	44.9	47.9	54.6	56.1	55.6
	724	44.8	48.4	54.4	54.3	52.6
VGG16	1024	55.2	57.3	64.5	66.9	67.4
	724	52.2	54.8	58.0	60.9	60.3

8 EXPERIMENTS

This section presents the results of our compact representation for image retrieval, evaluate the localization accuracy AML, and finally employ it for retrieval re-ranking.

Experimental setup. We evaluate the proposed methods on Oxford Buildings dataset (Philbin et al., 2007) and Paris dataset (Philbin et al., 2008), which are composed of 5063 and 6412 images, respectively. We refer to these datasets as Oxford5k and Paris6k. We additionally use 100k Flickr images (Philbin et al., 2007) to compose Oxford105k and Paris106k, respectively. A distractor set of 1 million images from Flickr (Jégou et al., 2010) is additionally used to go at larger scale. Retrieval performance is measured in terms of mean Average Precision (mAP). We follow the standard protocol and use the bounding boxes defined on the query images¹. These bounding boxes are also employed to evaluate localization accuracy. PCA is learned on Paris6k when testing on Oxford5k and vice versa. In order to be fair, we directly compare our results only to previous methods that do not perform learning on the test set.

The focus of our work is not to train a CNN, but to extract visual descriptors from its convolutional layers. We use networks widely used in the literature: AlexNet by Krizhevsky et al. (2012) and the very deep network (VGG16) by Simonyan & Zisserman (2014). We choose VGG16 instead of VGG19 because we observe that the latter does not always attain better performance while it has higher feature extraction cost. Our representation is extracted from the last pooling layer, which has 256 feature channels for AlexNet and 512 for VGG16. MatConvNet (Vedaldi & Lenc, 2014) is used to extract the features.

Localization accuracy. To evaluate the accuracy of AML, we employ pairs of Oxford5k query images and their corresponding positive images. We first perform exhaustive search to detect the globally optimal window. Then, we apply our speeded-up detector that evaluates fewer regions and in the end refines the best one. In both cases the approximate max-pooling is used for each window evaluation. We report Intersection over Union (IoU) with the optimal window and the percentage of windows evaluated compared to the exhaustive case. Results are shown in Table 1 (left). We provide a large speed-up while maintaining a high overlap with the optimal detection. Recall that our purpose is to apply this detector for fast re-ranking. Measuring IoU provides evidence for localization accuracy, however we observed that it does not directly impact retrieval performance. We finally set $s = 1.1$ and $t = 3$ for re-ranking usage.

In order to evaluate the localization accuracy with respect to ground-truth annotation we cross-match all 5 query images that exist per building. One of them is used as a query (cropped bounding box), while for the other we compare the detected region to the ground-truth annotation. Exhaustive evaluation achieves an IoU equal to 52.6% (52.9%) and the speeded-up approach achieves 51.3% (51.4%) on Oxford5k (Paris6k) datasets. The accuracy loss is limited, while the localization is approximately 180 times faster. AML provides a rough localization at low computational cost. Such

¹The query regions are cropped and then used as input to the CNN.

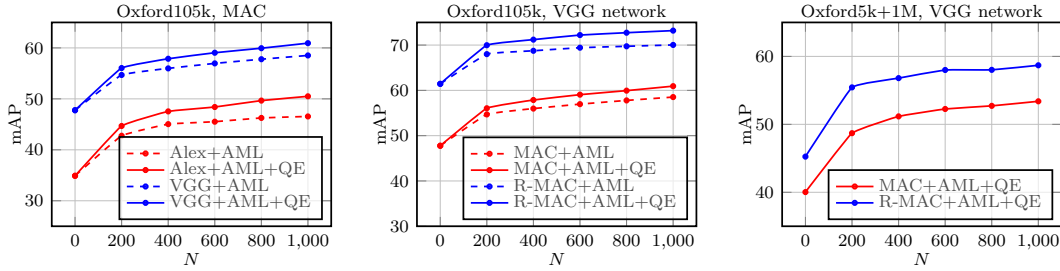


Figure 4: Performance of retrieval with re-ranking by AML versus number of re-ranked images on Oxford105k and Oxford5k combined with 1M distractor images.

Table 2: Performance comparison with state of the art. We report results for compact vector representations (left) and for retrieval approaches employing geometry, re-ranking, query expansion, or vector approximations (right). D = dimensionality. Our approaches are identified with bullets •.

Method	D	Oxf5k	Par6k	Oxf105k	Par106k	Method	Oxf5k	Par6k	Oxf105k	Par106k
Jégou & Zisserman (2014)	1024	56.0	-	50.2	-	Chum et al. (2011)	82.7	80.5	76.7	71.0
Jégou & Zisserman (2014)	128	43.3	-	35.3	-	Danfeng et al. (2011)	81.4	80.3	76.7	-
Babenko et al. (2014)	128	55.7	-	52.3	-	Mikulik et al. (2013)	84.9	82.4	79.5	77.3
Razavian et al. (2014b)	256	53.3	67.0	48.9	-	Shen et al. (2014)	75.2	74.1	72.9	-
Babenko & Lempitsky (2015)	256	53.1	-	50.1	-	Tao et al. (2014)	77.8	-	-	-
R-MAC •	256	56.1	72.9	47.0	60.1	Tolias et al. (2015)	80.4	77.0	75.0	-
R-MAC •	512	66.9	83.0	61.6	75.7	R-MAC+AML+QE •	77.3	86.5	73.2	79.8

a setup results in an average re-ranking query time of 2.9 sec using AlexNet, when re-ranking 1000 images with a single threaded implementation.

Retrieval and re-ranking. We evaluate retrieval performance using MAC and R-MAC compact representations. The MAC vectors are ℓ_2 -normalized, PCA-whitened and ℓ_2 -normalized once more, while the corresponding processing of the R-MAC is as described in Section 4. Table 1 (right) presents the results on Oxford5k. We evaluate different input image resolutions and observe that the original image size (1024) provides higher performance. Note that MAC is similar to the one proposed by Azizpour et al. (2014), however their process remains constrained by standard input size and aspect ratio. The proposed R-MAC gives a large performance improvement at no extra cost, as both feature vectors have exactly the same dimensionality. Regions of different scales are aggregated together, meaning that $L = 3$ combines regions at scales $l = 1$, $l = 2$, and $l = 3$. We set $L = 3$ in the following. In order to decompose the components of R-MAC, we construct R-MAC by aggregating only regions of $l = 3$. It achieves mAP equal to 63.0 on Oxford5k with VGG16. Aggregating both regions of $l = 2$ and $l = 3$ improves to 65.4. Finally, adding $l = 1$ (original R-MAC) performs 66.9 (see Table 1 right). Filtering time is 12 ms on average for Oxford105k.

Next, we employ AML for image re-ranking and conduct performance evaluation on Oxford105k by re-ranking up to 1000 images. The performance is consistently improved as shown in Figure 4. R-MAC brings a larger benefit and VGG16 performs better than AlexNet. Query expansion, as described in Section 6, improves the performance at low extra cost, since similarity is re-computed only for the re-ranked short-list. Finally, we carry out experiment at larger scale with 1M distractor images and present results in Figure 4. AML improves the performance by 13% mAP.

Examples of ranking using MAC and re-ranking using AML are presented in Figure 5. Recall that we only provide a rough object localization, since our main goal is to obtain improved image similarity. Furthermore, the provided localization is accurate enough for re-ranking.

Comparison to the state of the art. We compare the proposed methods to state-of-the-art performance of compact representations and approaches based on local features that perform precise



Figure 5: Examples of top retrieved images before (top) and after (bottom) re-ranking with AML. On the left we show the query image and depict the bounding box in blue color. When re-ranking is used, we present the top ranked images and report for each image its initial and final ranking. The localization window is shown in magenta, while positive/negative/junk images are depicted with green/red/yellow border.

descriptor matching, re-ranking or query expansion. Results are shown in Table 2². AlexNex and VGG16 are used to produce the 256D and 512D vectors for R-MAC, respectively. Regarding the compact representations, our short-sized R-MAC outperforms all other approaches. The better performance on Paris is inherited by the nature of the pre-trained networks; the baseline MAC with VGG achieves 55.2 on Oxford5k and 74.7 on Paris6k.

Unlike previous description schemes derived from CNN layers, our approach compete with the best approaches based on local features for geometric matching and query expansion. Our AML can even outperform them: while our results are lower on Oxford, we achieve the best performance on Paris and, to the best of our knowledge, outperform all published results on this benchmark. Higher scores on Paris6k are reported by Arandjelovic & Zisserman (2012) (91.0) and by Zhong et al. (2015) (91.5). These are achieved by learning the codebook on Paris6k itself and by performing pre-processing of the indexed dataset.

Discussion about other CNN-based approaches. Razavian et al. (2014b) propose to perform region cross-matching and accumulate the maximum similarity per query region. We evaluate this cross-matching process on the collection of regional vectors used in R-MAC; we simply skip the final aggregation process and keep the regional vectors individually. The cross-matching achieves 75.2% mAP on Oxford5k as a filtering stage, while re-ranking with AML on top of this acts in a complementary way and increases the performance up to 78.1%. However, cross-matching has

²Small differences of scores compared to the first version of the manuscript on arxiv are due to a slightly different evaluation protocol used before. Now, the evaluation protocol is the standard one for these datasets.

two drawbacks. Firstly, the region vectors have to be stored individually and increase the memory requirements by a factor of $|R|$, where $|R|$ is the number of extracted regions. Secondly, the complexity cost is linear in the number of indexed images and quite high since it requires to compute $|R|^2$ (e.g. 1024 (Razavian et al., 2014b)) inner products per image. The work of Razavian et al. (2014b) follows a non-standard evaluation protocol by enlarging the provided query bounding boxes. In addition, the cost of their feature extraction is extremely high since they feed 32 images of resolution 576×576 to the CNN. The recent work of Xie et al. (2015) is quite similar to theirs and is applied on both retrieval and classification.

Babenko & Lempitsky (2015) show that global sum-pooling on convolutional layer activations is better than max-pooling when the final image vectors are PCA-whitened. When whitening is not employed, then the latter is better. In the context of object localization we efficiently evaluate a large number of candidate regions on query time with AML. Performing whitening on each candidate region vector significantly increases the cost and is prohibitive for this task. We switch max-pooling to sum-pooling for both our proposed R-MAC and AML and test performance. Note that sum-pooling is a special case of our integral max-pooling with $\alpha = 1$. Switching to sum-pooling makes R-MAC perform 69.8 and R-MAC + AML + QE perform 76.9 on Paris106k. These scores are directly comparable to our scores in Table 2 and reveal that our choice is consistently better in all cases within our pipeline.

9 CONCLUSIONS

In this work, we re-visit both filtering and re-ranking retrieval stages by employing CNN activations of convolutional layers. Our compact vector representation encodes several image regions with simple aggregation method and is shown to outperform state-of-the-art competitors. Our localization increases the performance of the retrieval system that is initially based on a compact representation. The same CNN information adopted during the filtering stage is employed for re-ranking as well. Our approach competes with state-of-the-art methods that employ costly geometric matching or query expansion and we achieve the highest performance on Paris dataset, and provided a much better performance than existing approaches built upon CNN features. A very recent work (Arandjelovic et al., 2015) shows how MAC performance is improved by end-to-end fine tuning where the objective is based on MAC similarity.

REFERENCES

- Agrawal, Pulkrit, Girshick, Ross, and Malik, Jitendra. Analyzing the performance of multilayer neural networks for object recognition. In *ECCV*, 2014.
- An, Senjian, Peursum, Patrick, Liu, Wanquan, and Venkatesh, Svetha. Efficient algorithms for subwindow search in object detection and localization. In *CVPR*, 2009.
- Arandjelovic, Relja and Zisserman, Andrew. Three things everyone should know to improve object retrieval. In *CVPR*, Jun. 2012.
- Arandjelovic, Relja and Zisserman, Andrew. All about VLAD. In *CVPR*, Jun. 2013.
- Arandjelovic, Relja, Gronat, Petr, Torii, Akihiko, Pajdla, Tomas, and Sivic, Josef. Netvlad: Cnn architecture for weakly supervised place recognition. In *arXiv*, 2015.
- Avrithis, Yannis and Kalantidis, Yannis. Approximate gaussian mixtures for large scale vocabularies. In *ECCV*, 2012.
- Avrithis, Yannis and Tzouros, Giorgos. Hough pyramid matching: Speeded-up geometry re-ranking for large scale image retrieval. *IJCV*, 2014.
- Azizpour, Hossein, Razavian, Ali Sharif, Sullivan, Josephine, Maki, Atsuto, and Carlsson, Stefan. From generic to specific deep representations for visual recognition. In *arXiv*, 2014.
- Babenko, Artem and Lempitsky, Victor. Aggregating deep convolutional features for image retrieval. In *ICCV*, 2015.

- Babenko, Artem, Slesarev, Anton, Chigorin, Alexandr, and Lempitsky, Victor. Neural codes for image retrieval. In *ECCV*, Sep. 2014.
- Bentley, Joe. *Programming Pearls, 2/E*. Addison-Wesley Professional, 1999.
- Chen, Qiang, Song, Zheng, Feris, Rogerio, Datta, Amitava, Cao, Liangliang, Huang, Zhongyang, and Yan, Shuicheng. Efficient maximum appearance search for large-scale object detection. In *CVPR*, 2013.
- Chum, Ondrej, Mikulik, A., Perdoch, M., and Matas, J. Total recall II: Query expansion revisited. In *CVPR*, Jun. 2011.
- Danfeng, Qin, Gammeter, S., Bossard, L., Quack, T., and Gool, L. Van. Hello neighbor: Accurate object retrieval with k-reciprocal nearest neighbors. In *CVPR*, 2011.
- Dollár, Piotr, Tu, Zhuowen, Perona, Pietro, and Belongie, Serge. Integral channel features. In *BMVC*, 2009.
- Donahue, Jeff, Jia, Yangqing, Vinyals, Oriol, Hoffman, Judy, Zhang, Ning, Tzeng, Eric, and Darrell, Trevor. Decaf: A deep convolutional activation feature for generic visual recognition. In *arXiv*, 2013.
- Girshick, Ross. Fast r-cnn. In *arXiv*, 2015.
- Girshick, Ross, Donahue, Jeff, Darrell, Trevor, and Malik, Jitendra. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014.
- Gong, Yunchao, Wang, Liwei, Guo, Ruiqi, and Lazebnik, Svetlana. Multi-scale orderless pooling of deep convolutional activation features. In *ECCV*, 2014.
- Iandola, Forrest, Moskewicz, Matt, Karayev, Sergey, Girshick, Ross, Darrell, Trevor, and Keutzer, Kurt. Densenet: Implementing efficient convnet descriptor pyramids. In *arxiv*, 2014.
- Jégou, Hervé and Chum, Ondrej. Negative evidences and co-occurrences in image retrieval: The benefit of PCA and whitening. In *ECCV*, Oct. 2012.
- Jégou, Hervé and Zisserman, Andrew. Triangulation embedding and democratic aggregation for image search. In *CVPR*, 2014.
- Jégou, Hervé, Douze, Matthijs, and Schmid, Cordelia. Improving bag-of-features for large scale image search. *IJCV*, 87(3), Feb. 2010.
- Jégou, Hervé, Perronnin, Florent, Douze, Matthijs, Sánchez, Jorge, Pérez, Patrick, and Schmid, Cordelia. Aggregating local descriptors into compact codes. *Trans. PAMI*, Sep. 2012.
- Kalantidis, Yannis, Mellina, Clayton, and Osindero, Simon. Cross-dimensional weighting for aggregated deep convolutional features. In *arXiv preprint arXiv:1512.04065*, 2015.
- Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- Lampert, Christoph H. Detecting objects in large image collections and videos by efficient subimage retrieval. In *ICCV*, 2009.
- Lampert, Christoph H, Blaschko, Matthew B, and Hofmann, Thomas. Efficient subwindow search: A branch and bound framework for object localization. *Trans. PAMI*, 31(12):2129–2142, 2009.
- Lin, Zhe and Brandt, Jonathan. A local bag-of-features model for large-scale object retrieval. In *ECCV*, 2010.
- Lowe, David. Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110, Nov. 2004.
- Mikulik, Andrej, Perdoch, Michal, Chum, Ondřej, and Matas, Jiří. Learning vocabularies over a fine quantization. *IJCV*, 103(1), 2013.

- Oquab, Maxime, Bottou, Leon, Laptev, Ivan, and Sivic, Josef. Learning and transferring mid-level image representations using convolutional neural networks. In *CVPR*, 2014.
- Papandreou, George, Kokkinos, Iasonas, and Savalle, Pierre-André. Untangling local and global deformations in deep convolutional networks for image classification and sliding window detection. In *arXiv*, 2014.
- Philbin, James, Chum, Ondrej, Isard, Michael, Sivic, Josef, and Zisserman, Andrew. Object retrieval with large vocabularies and fast spatial matching. In *CVPR*, Jun. 2007.
- Philbin, James, Chum, Ondrej, Isard, Michael, Sivic, Josef, and Zisserman, Andrew. Lost in quantization: Improving particular object retrieval in large scale image databases. In *CVPR*, Jun. 2008.
- Radenović, Filip, Jegou, Herve, and Chum, Ondrej. Multiple measurements and joint dimensionality reduction for large scale image search with short vectors. In *ICMR*, 2015.
- Razavian, Ali Sharif, Azizpour, Hossein, Sullivan, Josephine, and Carlsson, Stefan. CNN features off-the-shelf: An astounding baseline for recognition. In *CVPRW*, 2014a.
- Razavian, Ali Sharif, Sullivan, Josephine, Maki, Atsuto, and Carlsson, Stefan. A baseline for visual instance retrieval with deep convolutional networks. In *arXiv*, 2014b.
- Ren, Shaoqing, He, Kaiming, Girshick, Ross, and Sun, Jian. Faster r-cnn: Towards real-time object detection with region proposal networks. In *arXiv*, 2015.
- Shen, Xiaohui, Lin, Zhe, Brandt, Jonathan, and Wu, Ying. Spatially-constrained similarity measure for large-scale object retrieval. *Trans. PAMI*, 36(6):1229–1241, 2014.
- Sicre, Ronan and Jurie, Frdric. Discriminative part model for visual recognition. *CVIU*, 141:28 – 37, 2015. ISSN 1077-3142.
- Simonyan, Karen and Zisserman, Andrew. Very deep convolutional networks for large-scale image recognition. In *arXiv*, 2014.
- Sivic, Josef and Zisserman, Andrew. Video Google: A text retrieval approach to object matching in videos. In *ICCV*, 2003.
- Tao, Ran, Gavves, Efstratios, Snoek, Cees GM, and Smeulders, Arnold WM. Locality in generic instance search from one example. In *CVPR*, 2014.
- Tolias, Giorgos, Avrithis, Yannis, and Jégou, Hervé. Image search with selective match kernels: aggregation across single and multiple images. *IJCV*, 2015.
- Uijlings, Jasper, Van de Sande, Koen, Gevers, Theo, and Smeulders, Arnold. Selective search for object recognition. *IJCV*, 104(2):154–171, 2013.
- Van de Sande, Koen EA, Snoek, Cees GM, and Smeulders, Arnold WM. Fisher and VLAD with flair. In *CVPR*, 2014.
- Vedaldi, Andrea and Lenc, Karel. Matconvnet-convolutional neural networks for matlab. In *arXiv*, 2014.
- Viola, Paul and Jones, Michael. Robust real-time object detection. *IJCV*, 4:34–47, 2001.
- Xie, Lingxi, Tian, Q, Hong, R, and Zhang, B. Image classification and retrieval are one. In *ICMR*, 2015.
- Zhong, Zhiyuan, Zhu, Jianke, and Hoi, Steven CH. Fast object retrieval using direct spatial matching. *IEEE Trans. on Multimedia*, 17(8):1391–1397, 2015.